

VIRTUALIZATION OF iSCSI STORAGE

FIELD OF THE INVENTION

This invention is directed to the field of IP based storage networks. It is more particularly directed to the virtual access of iSCSI (Internet Protocol -- Small Computer Systems Interconnect) storage devices.

BACKGROUND OF THE INVENTION

Storage-area networks, or SANs, are gaining in popularity because they promise to curb the rising costs of storage management by enabling wider sharing of storage devices and the consolidation of storage resources under centralized administrative control. The promise of storage-area networks to simplify management relies on their ability to virtualize storage devices, separating the *virtual* or logical view of storage from the *physical* view. Storage virtualization allows administrators to deal and manage the simpler *virtual view*, while the storage management system handles the complexities of how that view is implemented on top of physical resources. Therefore, a high-performance and secure storage virtualization solution is crucial for such storage networks.

When storage virtualization is employed, the applications, which in this context refer to the file servers and database servers and any other application accessing block-level devices, are presented with a virtual storage space which has the required performance and availability requirements. The implementation and management of storage to provide the requisite levels of performance and

1 availability is hidden and can change underneath the covers without application
2 knowledge or participation.

3 Virtual storage provides the illusion of expandable storage space thereby isolating
4 the clients from the management of physical storage resources, such as disks, disk
5 arrays and tapes. While the underlying physical devices have fixed and limited
6 capacity, a virtual storage repository can expand its capacity on a per need basis,
7 and can improve its performance by changing the underlying physical storage
8 devices used. Another advantage of virtualization is that it allows for load
9 balancing to occur without host participation. When the physical blocks are be
10 moved to balance load, but application-visible names do not have to be changed.
11 Furthermore, storage virtualization allows for the view (namespace) of visible
12 storage to be customized on a per-host basis and security and access control
13 policies to be managed on a per-host basis.

14 The basic idea of storage virtualization is to provide a layer of indirection,
15 mapping virtual storage blocks to physical blocks. This invention concerns
16 storage-area networks which use iSCSI devices. iSCSI is an TCP/IP based
17 protocol to carry SCSI commands over an IP network between hosts and storage
18 devices. . Furthermore, we suppose that the SCSI storage devices are connected
19 via a switched SAN within a data center. SAN gateways are placed at the edge of
20 the SAN to provide the virtual storage abstraction to applications running on the
21 hosts. All traffic to the devices goes through one of the SAN gateways.

22 In such a system, a good virtualization solution should achieve the following
23 goals:

- 24 • Security and access control: Security is critical to protect storage.
- 25 • High-performance: Avoiding data copies and connection management at
26 the virtualization gateway increases bandwidth.

- 1 • Manageability of storage: Security protocol upgrades, storage migration
2 should be easy to do.

3 **SUMMARY OF THE INVENTION**

4 It is thus an aspect of the present invention to divide each virtual logical unit
5 (LUN) into block ranges of fixed size, with each range mapped on to a physical
6 LUN on a single device.

7 It is another aspect of the invention to export to each host a unique IP address for
8 a given virtual LUN. The host accesses different block ranges within the virtual
9 LUN via different TCP port numbers but via the virtual LUN's IP address.

10 Still another aspect of this invention is to use a gateway to perform access control
11 and a level of virtualization by mapping virtual (IP, port #) pairs in IP packets sent
12 by the host onto actual (IP, port#) pairs of physical storage devices.

13 Other aspects and a better understanding of the invention may be realized by
14 referring to the detailed description.

15 **BRIEF DESCRIPTION OF THE DRAWINGS**

16 These and other aspects, features, and advantages of the present invention will
17 become apparent upon further consideration of the following detailed description
18 of the invention when read in conjunction with the drawing figures, in which:

1 Fig.1 describes a storage area network, in which storage devices are connected via
2 a network to end hosts through a storage area network gateway;

3 Fig. 2 represents an enhanced iSCSI storage device with multiple LUNs along
4 with a mapping of TCP port numbers to LUNs, with each mapping associated
5 with access rights, in accordance with the present invention;

6 Fig. 3 shows the enhancements required at the iSCSI layer on the host, in
7 accordance with the present invention;

8 Fig. 4a shows the use of a router with address translation capability as a storage
9 virtualization gateway, in accordance with the present invention;

10 Fig. 4b shows the use of address translation tables at the storage virtualization
11 gateway, in accordance with the present invention;

12 Fig. 5a shows the use of a router with address translation and IPSec processing
13 capabilities as a secure storage virtualization gateway, in accordance with the
14 present invention;

15 Fig. 5b shows the different packet processing capabilities supported by the secure
16 virtualization gateway, in accordance with the present invention;

17 Fig. 6 shows migration of storage blocks between two physical storage devices
18 and changes in the address translation table at the virtualization gateway such that
19 the host remains unaffected by this migration, in accordance with the present
20 invention, in accordance with the present invention;

1 Fig. 7a shows the virtualization support modules at the host, in accordance with
2 the present invention;

3 Fig. 7b shows the virtualization support modules at the gateway, in accordance
4 with the present invention; and

5 Fig. 7c shows the virtualization support modules at the storage device, in
6 accordance with the present invention.

7 **DESCRIPTION OF THE INVENTION**

8 Figure1 shows a storage area network (SAN) with virtualization gateways. A
9 storage area network (SAN) is composed of storage devices (104, 105), gateway
10 (106) and hosts (101,102,103). Gateways are on the edge of the SAN. Hosts talk
11 iSCSI to the gateway. Gateways talk iSCSI to the devices. In such a system, hosts
12 acts as clients requesting data blocks, devices as block servers. Gateways perform
13 functions such as virtualization and access control. A SCSI (iSCSI) command
14 addresses a logical unit number (LUN), specifies an offset and the number of
15 blocks, to read and write including the starting block. When virtualization is
16 used, the arguments specified by the host in the SCSI command are actually
17 virtual. They need to be mapped to their physical counterparts. In this invention,
18 the term LUN will be used to refer to the logical unit itself, as well as to the
19 identifier for the logical unit, [i.e. the logical unit number] as used by those skilled
20 in the art.

21 The gateways fulfill three functions, the first and primary function is routing. The
22 gateways are commodity network switches or routers. The second function is

1 assisting with translations (to support storage virtualization). The third function is
2 ensuring proper access control and security at the edge of the network so that the
3 devices do not have to implement a sophisticated authentication or security
4 protocols. The number of gateways is expected to be smaller than the number of
5 devices and therefore more manageable. Constraining security functions to the
6 gateways reduces cost by limiting the nodes where secret keys are stored and
7 where cryptographic accelerators are added, simplifies the devices and the
8 management or update of security protocols.

9 A straightforward implementation of a virtualization gateway for iSCSI devices
10 and hosts is to terminate TCP connections from the host, retrieving the SCSI
11 command from the host packets. The gateways can then translate the virtual
12 access to a physical access and use one or more TCP connections to the physical
13 devices to transmit the modified physical commands, then merge and return the
14 results to the host. This of course requires data copying, connection management
15 and full processing through the TCP/iSCSI and SCSI stacks at the gateway.
16 Consequently, this load limits the performance (throughput) of the gateway.

17 Our solution relies on limited support performed at the host and some checks and
18 network address translations at the gateway to achieve direct access with little
19 connection management and no data touching at the gateway. To allow the
20 gateway to perform the routing and access checks without parsing the SCSI
21 command inside the packet, the gateway uses the following scheme. The gateway
22 uses the port numbers publicized to the host, and which the host uses in every
23 subsequent packet to decode the target physical logical unit number (LUN)
24 identifier the packet should be routed to.

1 The gateway publicizes tables containing metadata about each virtual LUN to the
2 host. These tables specify a different port for each block range within the virtual
3 LUN. Each such range is mapped onto a different physical LUN. Multiple
4 physical LUNs may reside on the same physical device but they are associated
5 with different ports and can be migrated to other devices independently of each
6 other. As a result, migrations and reconfiguration will not require host
7 notification. Only the maps used by the gateway need to be updated. When
8 receiving a packet that is part a TCP connection to a particular block range, all
9 the gateway has to do is steer it to the proper physical LUN by rewriting the IP
10 and port numbers in the packet headers. The gateway then translates an incoming
11 packet header <src address, virtual dest addr, gateway-fake port number> to <src
12 addr, physical device IP addr, physical device port number> where the dest addr is
13 a function of source address, virtual dest addr and dest port number. The
14 virtualization gateway is thus provided by a regular network address translation
15 (NAT) box.

16 As shown in Figure 2, a storage device supports multiple physical LUNs with a
17 different TCP port number associated with each physical LUN. One aspect of the
18 invention is that all iSCSI commands received on a given TCP port of a storage
19 device correspond implicitly to the physical LUN associated with that port, and
20 while the offset and block numbers in the iSCSI command are significant, the
21 LUN identifier in the command is ignored. Figure 2 shows a storage device 201
22 which supports physical logical units LUN0 (207), LUN1 (208) and LUN1 (209)
23 which received iSCSI commands on TCP port numbers port0 (204), port1(205)
24 and port2 (206) respectively. The storage device is connected to a virtualization
25 gateway through a communication link 203. The table 202 stores access rights for
26 each physical LUN.

Figure 3 shows the steps performed by the host to process a SCSI command request. The host caches a table 301 which associates a single IP address for each virtual LUN, and the SCSI command parameters (LUN, Starting Block, Number of Blocks), shown as item 309 in the figure, are translated by the host to one or more iSCSI commands (Physical LUN, Remapped Starting Block, Remapped #Blocks) on one or more TCP connections, all to the same IP address, but different port numbers, with each iSCSI connection corresponding to a different TCP port number. In this figure, the table shows two entries 302 and 303, corresponding to VLUN#0 and VLUN#1, which are mapped to virtual IP addresses IP0 and IP1, respectively. Each entry maps block ranges within a VLUN to specific TCP port numbers. Commands issued by the SCSI layer 305 at the host, such as 309 in Figure 3, are translated by the enhanced iSCSI layer 306 by looking up the appropriate entry in the table 301. The packets are then handed over to the TCP/IP layer 307 at the host, followed by an optional IPsec layer 308 which is responsible for setting up a secure tunnel with the virtualization gateway, as will be discussed in Figure 5.

The invention requires that a device having multiple physical LUNs associate a port with each LUN. All commands received on a port are assumed implicitly to target the corresponding LUN associated with that port. Thus, Note that the commands issued by the host even when split into multiple commands for different chunks (different physical LUNs) will have the VLUN identifier in the command arguments embedded in the SCSI command within the TCP packet.

Once the host-side command rewriting is performed, outgoing SCSI commands use the correct offsets within the physical LUNs. The command is sent to the gateway, and the gateway routes the packet to the proper physical device on which the physical LUN onto which the chunk is mapped resides. As shown in

Figure 4, the gateway 402 performs an IP header rewriting of the destination IP address and port number, without touching the data or terminating TCP connections. The gateway indexes into a local table 401 to retrieve the address, port translations. If an mapping is absent, then the host was not allocated that address and the gateway drops the packet. This allows the gateway to enforce access control such that a host can access only the address space that has been exported to it. The table 401 maps <Virtual IP address, TCP port> on packets incoming from the hosts to <IP address, TCP port> corresponding to the physical LUN of the physical storage devices.

The gateway uses the standard IPSEC protocol to ensure authenticated optionally encrypted and private traffic between itself and the host. Also the gateway performs authorization checks. It verifies that a command to a target physical unit is from a host that is authorized to issue such a command. This is achieved as follows. The gateway has a map providing what physical logical units are accessible to what hosts. Upon receiving an authenticated IP packet from a host, it performs a quick lookup in a hash-table indexed by (src-ip, port #) to retrieve the rights of the host with source ip address src-ip to the physical logical unit uniquely identified with gateway-port#. If an entry exists providing the host the write to access the command, the packet is forwarded, simply translating the IP address field in the packet to the IP address of the physical device and changing the port# from gateway-port# to the recorded port number of the physical logical unit.

Through IPsec, we can support different levels of security, simple authentication, authentication plus integrity of packet (thereby ensuring command & data integrity) or full privacy (through payload encryption). Note that the devices need not have any IPsec or encryption support. Thus, they do not need to be upgraded

1 whenever a weakness in the protocol or encryption method is detected. All
2 security work is restricted to the much fewer gateways.

3 One advantage of storage virtualization is that storage managers, servers that are
4 deployed within the SAN to move and reconfigure storage to balance load and
5 capacity across devices, can do so without host coordination, involvement or
6 support. Therefore, any virtualization solution must support the on-line
7 reconfiguration of storage. The problem with storage migration tasks is that they
8 move data blocks around and therefore the maps that translate a virtual block-id to
9 a physical block-id must be updated to reflect the new location of a physical block
10 that has been recently moved.

11 Figure 6 shows how the above storage virtualization scheme is used to migrate
12 logical units between storage devices without requiring the host to participate in
13 the migration process. The host 604 has a virtualization map, which maps the
14 accesses to different blocks of VLUN#0 to different TCP port numbers on IP
15 address IP_v0, as shown in 605. In this example, VLUN#0 is shown to contain
16 1000 blocks, all of which are mapped to port0. This is initially mapped to LUN0
17 of storage device 606 with a physical IP address IP1; commands for LUN0 are
18 received on port0 on IP1. The virtualization gateway 608 initially translates
19 packets from the host with source IP address IP0, according to entry 602 in its
20 translation table 601. The virtual destination IP address, IP_v0, is replaced by IP1
21 and the destination port number port0 is unchanged. This is because accesses to
22 the virtual storage device VLUN#0 by the host is mapped to the physical unit
23 LUN0 of storage device 602 with physical address IP1.

24 Now, lets assume that this mapping needs to be changed and the accesses to
25 VLUN#0 by host IP0 should be remapped to LUN2 of storage device 603 with IP

address IP2; LUN2 of storage device 603 receives SCSI commands on port2. To facilitate this remapping, the entry 602 at the gateway's translation map 601 is replaced by entry 603. Consequently, the destination address IP_v0 on incoming packets at the gateway is replaced by IP2, and the destination port number port0 is replaced by port2, and TCP/IP packets containing iSCSI data/commands that were earlier being sent to LUN0 (port0) of storage device IP1 are now being sent to LUN2 (port2) of storage device 607 without changing any entry of the translation map 605 at the host 604. Since iSCSI operates over TCP connections, the host will receive a TCP reset the first time it sends a packet to the storage device 607, since it is unaware of the migration, ie remapping of its virtual storage unit VLUN#0. As a result, the TCP connection will be automatically reset, i.e the existing connection will be torn down and a new connection will be set up with the same destination address IP_v0 (as far the host is concerned). SCSI commands/data can now be exchanged over this connection between the host and the storage device 607. Physical communication links between the host and gateway, and between the gateway and the two storage devices are shown as 609, 610 and 611.

Figure 7a shows the different modules implementing the invention at the host. A virtualization module (701) includes a control module (702) and a driver module (703). Figure 7b shows the address translation module at the gateway (704), while Figure 7c shows the conversion module (705) required at the storage device. These modules can be implemented in a manner known to those skilled in the art.

The present invention can be realized in hardware, software, or a combination of hardware and software. A visualization tool according to the present invention can be realized in a centralized fashion in one computer system, or in a distributed fashion where different elements are spread across several interconnected

1 computer systems. Any kind of computer system - or other apparatus adapted for
2 carrying out the methods and/or functions described herein - is suitable. A typical
3 combination of hardware and software could be a general purpose computer
4 system with a computer program that, when being loaded and executed, controls
5 the computer system such that it carries out the methods described herein. The
6 present invention can also be embedded in a computer program product, which
7 comprises all the features enabling the implementation of the methods described
8 herein, and which - when loaded in a computer system - is able to carry out these
9 methods.

10 Computer program means, or computer program, in the present context include
11 any expression, in any language, code or notation, of a set of instructions intended
12 to cause a system having an information processing capability to perform a
13 particular function either directly or after conversion to another language, code or
14 notation, and/or reproduction in a different material form.

15 Thus the invention includes an article of manufacture which comprises a
16 computer usable medium having computer readable program code means
17 embodied therein for causing a function described above. The computer readable
18 program code means in the article of manufacture comprises computer readable
19 program code means for causing a computer to effect the steps of a method of this
20 invention. Similarly, the present invention may be implemented as a computer
21 program product comprising a computer usable medium having computer
22 readable program code means embodied therein for causing a a function described
23 above. The computer readable program code means in the computer program
24 product comprising computer readable program code means for causing a
25 computer to effect one or more functions of this invention. Furthermore, the
26 present invention may be implemented as a program storage device readable by

2025 RELEASE UNDER E.O. 14176

1 machine, tangibly embodying a program of instructions executable by the
2 machine to perform method steps for causing one or more functions of this
3 invention.

4 It is noted that the foregoing has outlined some of the more pertinent objects and
5 embodiments of the present invention. This invention may be used for many
6 applications. Thus, although the description is made for particular arrangements
7 and methods, the intent and concept of the invention is suitable and applicable to
8 other arrangements and applications. It will be clear to those skilled in the art that
9 modifications to the disclosed embodiments can be effected without departing
10 from the spirit and scope of the invention. The described embodiments ought to
11 be construed to be merely illustrative of some of the more prominent features and
12 applications of the invention. Other beneficial results can be realized by applying
13 the disclosed invention in a different manner or modifying the invention in ways
14 known to those familiar with the art.